



Middlesex University Research Repository

An open access repository of
Middlesex University research

<http://eprints.mdx.ac.uk>

Gao, Xiaohong W. (2014) Feature-wise representation for both still and motion 3D medical images. In: 2014 IEEE Southwest Symposium on Image Analysis and Interpretation (SSIAI), April 6-9, 2014, San Diego, USA.

Available from Middlesex University's Research Repository at
<http://eprints.mdx.ac.uk/13301/>

Copyright:

Middlesex University Research Repository makes the University's research available electronically.

Copyright and moral rights to this thesis/research project are retained by the author and/or other copyright owners. The work is supplied on the understanding that any use for commercial gain is strictly forbidden. A copy may be downloaded for personal, non-commercial, research or study without prior permission and without charge. Any use of the thesis/research project for private study or research must be properly acknowledged with reference to the work's full bibliographic details.

This thesis/research project may not be reproduced in any format or medium, or extensive quotations taken from it, or its content changed in any way, without first obtaining permission in writing from the copyright holder(s).

If you believe that any material held in the repository infringes copyright law, please contact the Repository Team at Middlesex University via the following email address:

eprints@mdx.ac.uk

The item will be removed from the repository while any claim is being investigated.

Feature wise representation for both still and motion 3D medical images

Xiaohong Gao

School of Science and Technology,
Middlesex University, London, NW4 4BT, United Kingdom
x.gao@mdx.ac.uk

Abstract— *With the arrival of the state of the art medical imaging equipment, a plethora of images are acquired not only in higher dimensions (3D+) but also with various presenting forms of either still or motion, complicating data management systems even further. This paper offers, from an application point of view, representations of content features from both still 3D MR brain images and 3D ultrasound cardiac video sequences by demonstrating a developed online content-based image retrieval system, MIRAGE. The approaches of 3D SIFT coupled with sparse code have been appointed to facilitate the representation of image features, whereas widely applied four texture based approaches are also implemented to allow users' benefit of different retrieving intentions.*

Keywords—*Imaging, representing biomedical knowledge, multi-media image retrieval.*

I. INTRODUCTION

Three dimensional images, including both still pictures and motion videos, are at present a common form in hospitals in assisting clinicians performing diagnosis or intervention, such as 3D MR brain images or ultrasonic 2D video clips with time scale being the third dimension. With regard to still images, structured templates usually exist, containing constructive geometric properties, such as shape or texture, the characteristics that can be employed to standardise those images in terms of their content features. For motion videos that are of a function of both space and time, however, additional data might be needed to ensure that the starting time is in the same cycle for all the datasets in a database.

With regard to 3D still brain images, many methods for feature representation of content have been developed, including intensity-based [1], physiological-information-based [2], and textured-based [3]. With the application of the Talairach [4] brain atlas, intensity-based approach takes advantages of spatial references to concert a region-based retrieval, in which a regional or a volumetric data is expressed as $\langle x, y, z, \text{value} \rangle$. Since this representation is reminiscent of image matrix, depending on the size of each region/volume, the presentation volume of each image can be as large as the size of an image dataset itself, to a certain extent, missing the point of representing images using features to reduce redundant information and therein being prone to artefacts. On the other hand, by using physiological-information-based retrieval, a number of semantic contents can be dealt with by employing

physiological kinetic features of images [2]. Although effective, this method is very discipline-constrained and heavily relies on the additional supply of extra information, such as blood samples in order to derive plasma time activity (PTA) curves, gaining prospective in searching for alternative approaches.

Texture-based approach to extract features from 3D MR brain images is originally applied by the employment of 2D Gabor filter [5], which in essence remains a 2D approach. After further extension and tailoring, 3D textures are obtained from four well known approaches in [3], including Local Binary Pattern (LBP), Grey Level Co-occurrence Matrices (GLCM), Wavelet Transforms (WT) and Gabor Transforms (GT). For the application to 3D brain images using texture-based approaches, the prerequisite resides on the spatial normalisation to a standard template. In this way, the comparison of each sub-volume (cube) within the same location between images can be conducted under the assumption that geometrically, all brains bear similar anatomic structure. In practice, however, due to lesioned brains sustaining appreciable distortions, the texture features even from normalized images cannot be unique in relation to representations, which is also true for video images. To overcome this deficit, in this study, key point based feature representation by the appointment of scale invariant feature transformation (SIFT) coupled with a machine learning technique of sparse coding is implemented, towards developing an online system of content-based image retrieval (CBIR) system for medical images (<http://image.mdx.ac.uk/time/demo.php>).

SIFT descriptors [6], being invariant to geometric transformations of translation, scaling and rotation, provide robust feature matching mechanisms across a substantial range of changes in illumination while with the presence of noise. It is therefore widely applied in the domain of object recognition and image stitching. In addition, a 3D extension of the SIFT algorithm has recently been proposed in the literature on 3D volumetric data analysis, such as, action recognition in video volumes [7], object recognition in CT complex volume [8], 3D medical registration and panoramic medical image stitching [9]. Evidently, each type of dataset present their own unique set of characteristics that subsequently require the application of SIFT in a different way as to the extraction of potential features.

The remainder of the paper is structured as follows. Section II addresses mathematical formulae underling the methodology, which is then followed by the results on the application of CBIR for brain images and classification for cardio videos given in Section III. The paper is then concluded in Section IV.

II. METHODOLOGY

The procedures involve sparse coding of 3D salient content features and the creation of a codebook of visual dictionaries.

A. Pre-processing

Local visual feature selection usually remains the first task to conform. For an MR brain image, after spatially normalised to an MNI (Montreal Neurological Institute) template, it is spatially partitioned equally into non-overlapped sub-cubes of 512 ($=8 \times 8 \times 8$). In this way, the retrieval practice can be conducted by focusing on each corresponding sub-cube without the need of cross-comparison, saving considerable retrieving time.

On the other hand, for motion video clips, low-level interest-points are found out first as potential candidates for cross-image comparisons. To detect a spatial-temporal interest point, the technique of Cuboid detector [10] are opted for in an attempt to overcome the inability that many other methods suffer from with reference to incorporating temporal information.

With respect to still brain images, 3D SIFT descriptors are applied to enumerate local visual features by computing a 3D gradient orientation histogram for each of 512 sub-volumes. In doing so, further division is performed on each sub-volume to create eight ($=2 \times 2 \times 2$) sub-blocks, upon each of which, the magnitude and orientation of its gradient are calculated, by the application of 1D Haar wavelet transform in each of x, y and z directions respectively, forming bins of orientations accumulating the magnitude values that share the same gradient orientation. Subsequently, based on the tessellation technique, eighty bins of orientation are rendered in a 3D orientation sphere for each sub-block, leading to a feature vector of 640 dimensions ($=2 \times 2 \times 2 \times 80$) for each sub-volume.

Likewise, a 3D SIFT descriptor consisted of 640 elements is extracted as a feature vector in video images, in this case to represent each interest voxel instead of sub-volume. Centered at each interest voxel, detected by using the Cuboid detector, a $12 \times 12 \times 12$ volume of neighbourhood is selected and then divided into 8 ($=2 \times 2 \times 2$) sub-volumes that subsequently undertake the same procedure of calculation of gradients as explained above.

B. Visual Vocabulary Construction Using Sparse Coding

Once 3D SIFT features, i.e., the candidates for unit elements, or “words” in a visual dictionary, are accounted

for from each sub-volume or voxel, sparse coding follows to allow the creation of dictionary of visual ‘words’.

In theory, by modeling data vectors as a sparse linear combination of a set of basic elements or ‘words’, sparse coding [11] encodes each descriptor of an image by solving the optimization problem as formulated in Eq.(1).

$$\min_{U,V} \sum_{m=1}^M \|x_m - Vu_m\|_2^2 + \lambda \|u_m\|_1$$

$$\text{Subject to: } \|v_i\| \leq 1, \quad \forall i = 1, \dots, K \quad (1)$$

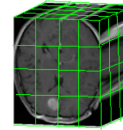
where $X = [x_1, x_2, \dots, x_M]$ ($x_m \in R^{dx1}$) refers to a set of 3D SIFT descriptors as described above from a 3D training dataset. $V = [v_1, v_2, \dots, v_K]$ ($v_i \in R^{dx1}$) indicates the K bases, also known as the dictionary or codebook; and $U = [u_1, u_2, \dots, u_M]$ ($u_m \in R^{K \times 1}$) denotes sparse codes for images based on codebook V . In addition, M refers to the number of total training samples in the training dataset, whereas λ indicates the constant coefficient that is generated by L_1 norm ($\|\cdot\|_1$) regularization, which acts as the penalty function to produce sparse coefficients and being robust to irrelevant features.

C. Max-pooling for Image Representation

To take into account of spatial location of local features, a pooling technique is applied to the SIFT sparse code that is calculated from different regions of an image to create the representation of an image by concatenation. In this study, max-pooling technique is employed as opposed to average-pooling by choosing the max value from a set of inputs as illustrated in Eq. (2).

$$z_i = \max\{u_{ij}, j = 1, 2, \dots, S\} \quad (2)$$

where u_{ij} in Eq.(2) indicates the i^{th} ($i \in [1, K]$) 3D SIFT sparse code for the j^{th} ($j \in [1, S]$) sub-volume or interest point within a pooling region, with K indicating the size of codebook V , whereas S refers to the total number of the sub-volumes in each pooling region. Figure 1 schematically demonstrates the pooling process, i.e., an MR cubic image being equally divided into 64 ($= 4 \times 4 \times 4$) regions, whilst 12 pooling regions are drawn up from a clip of ultrasonic videos to acknowledge the differences of the structure of a moving heart at three levels.



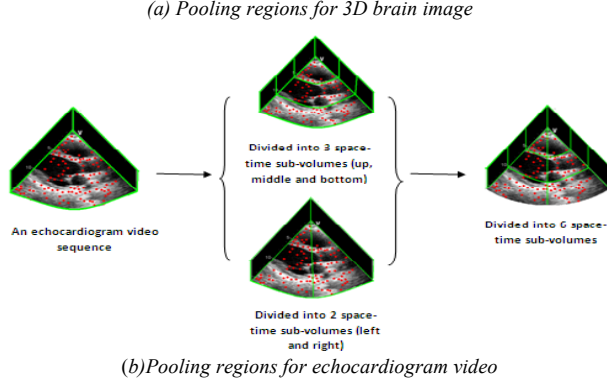


Figure 1. The pooling regions for 3D MR brain images (a) and echocardiograms (b)

D. Comparisons between Datasets

For still images, both histogram intersection and Chi-square histogram are applied to measure the degree of similarity between two images. Whereas for video images, the comparison is based on viewpoint classification which is performed using an approach of multiclass SVM (Support Vector Machine) with a linear kernel that is calculated in Eq. (3).

$$k(\bar{F}_i, \bar{F}_j) = \bar{F}_i^T \bar{F}_j \quad (3)$$

where \bar{F}_j is the feature representation of video j . With regard to binary classification, an SVM aims to learn a decision function based on the training dataset, which is defined in Eq. (4).

$$f(\bar{F}) = \sum_{i=1}^n a_i k(\bar{F}_i, \bar{F}) + b \quad (4)$$

In order to obtain an extension to a multi-class SVM, the trained videos are represented as $\left\{ \left(\bar{F}_i, l_i \right) \right\}_{i=1}^n$, where $l_i \in \{1, 2, \dots, L\}$ denotes the class label of trained video i . One-against-all strategy is applied to train the total number of L , the binary classifiers.

III. RESULTS

A. Dataset Collection

In this investigation, the data at our disposal are of both 3D MR still brain images and 3D (i.e., 2D video) echocardiographs, all in the format of DICOM. The characteristics of these datasets are given in Table 1.

TABLE 1. THE INFORMATION RELATED TO THE DATASETS AT THE DISPOSAL.

	Still Brain Images	Ultrasonic Motion Video Cardiac Clips
Normal	34	14
Abnormal	86	58

Subject	120	72
Resolution	500×500×45 (mm ³)	434 × 636 (pixel ²) × 26 (frames)
Imaging Tool	GE Genesis_Signa 15000T (1.5-T) whole-body MR	GE Vivid 7 Ultrasound
Total	120	219

Additionally, the ground truth is based on the locations of lesions from brain images and eight viewpoints of heart video clips, which are marked by clinicians.

B. Training

In the training stage, a 3D brain codebook composed of 64 sub-codebooks is obtained from 500 descriptors randomly selected from 960 sub-volumes.

With regard to echo-cardiac video images, 80,000 interest points firstly detected by Cuboid detector are randomly selected from the video clips (n=219) as a training dataset for the generation of the codebook.

C. Comparison Results

As discussed in the INTRODUCTION, 3D still images can also be represented using 3D texture-based approaches. Therefore comparison with the four popular texture approaches [3], is also carried out, with the comparison results given in Table 2, by which retrieval task focuses on the location of lesions by using the measure of mean average precision (MAP).

TABLE 2. MAP VALUES FOR THE FIVE 3D APPROACHES INCLUDING GLCM, WT, GT, LBP AND SIFT

Methods	Without sparse coding	With sparse coding	
		Histogram intersection	Chi-squared histogram
3D GLCM	0.3034	0.3291	0.3510
3D WT	0.3096	0.3375	0.3687
3D GT	0.3074	0.3863	0.3954
3D LBP	0.3308	0.4027	0.4012
3D SIFT	0.3959	0.4013	0.4098

As presented in Table 2, the approach of 3D SIFT that has been furthered in this research outperforms the other four with the average MAP value of 0.4098 according to Chi-squared histogram distance. In addition, the implementation of sparse coding improves the performance of all five approaches, specifically for 3D GT with MAP value increasing from 0.3074 to 0.3954, implying the significance of contribution of machine learning technique to the presentation of visual features.

For video images, the retrieval is based on the viewpoint. Therefore in essence the procedure is of a classification. Table 3 illustrates the confusion matrix of the classification based on eight standard viewpoints.

TABLE 3. CONFUSION MATRIX FOR 8 ECHOCARDIOGRAM VIEW CLASSIFICATION WHERE AR=ACCURACY RATE, ER=ERROR RATE, AND THE VERTICAL AXIS REFERS TO GROUND TRUTH

	Classification Results	AR
--	------------------------	----

